

# Lecture 2: Neural Networks

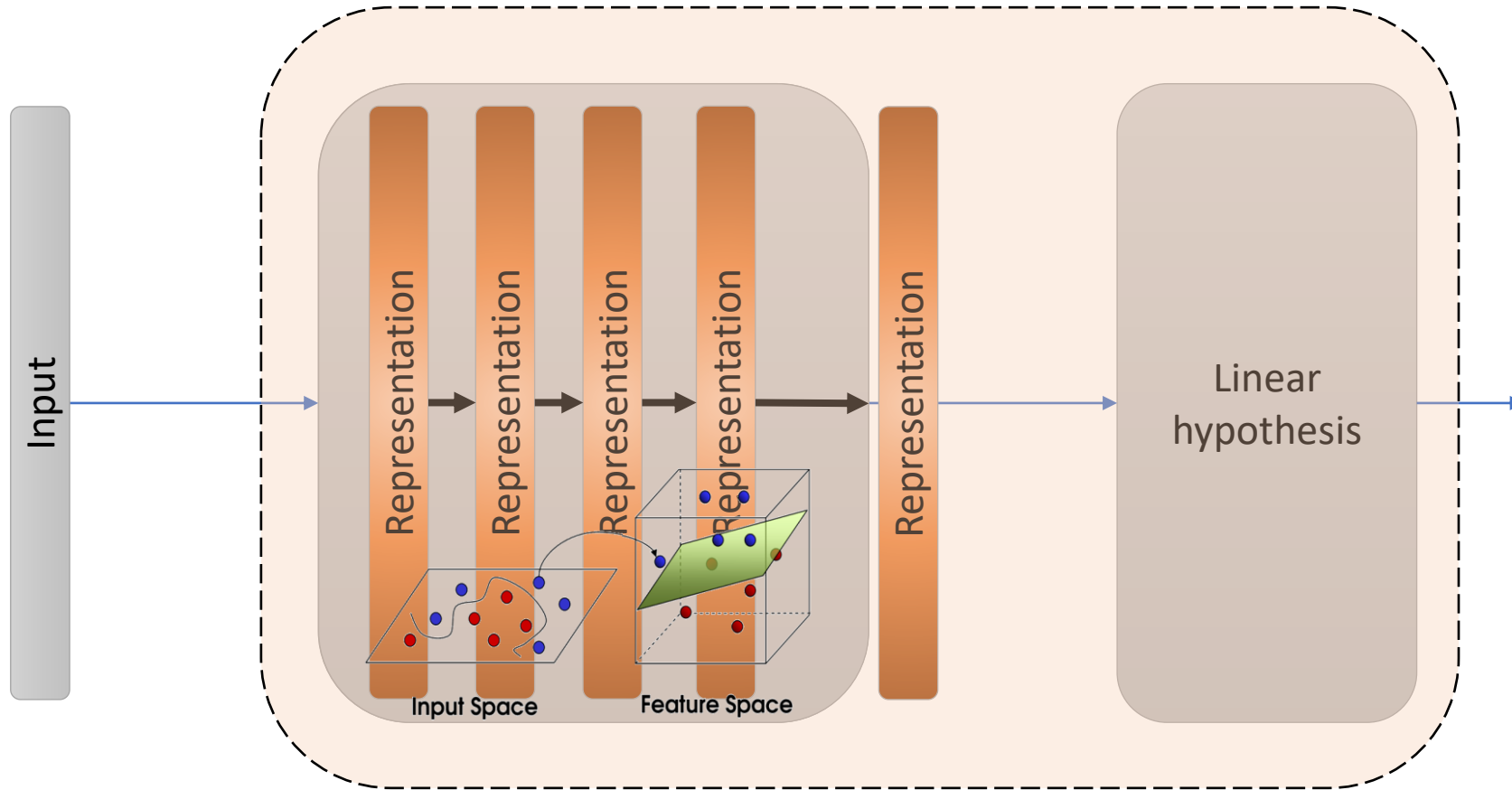


# Today:

- Revisit feature transform (5%)
- What is a neural net? (10%)
- Derivatives and chain-rule reminder (10%)
- Training a vanilla network (back-prop) (40%)
- Differential computational graph (25%)
- Demo (10%)



# Feature transform



Non-linear hypothesis!

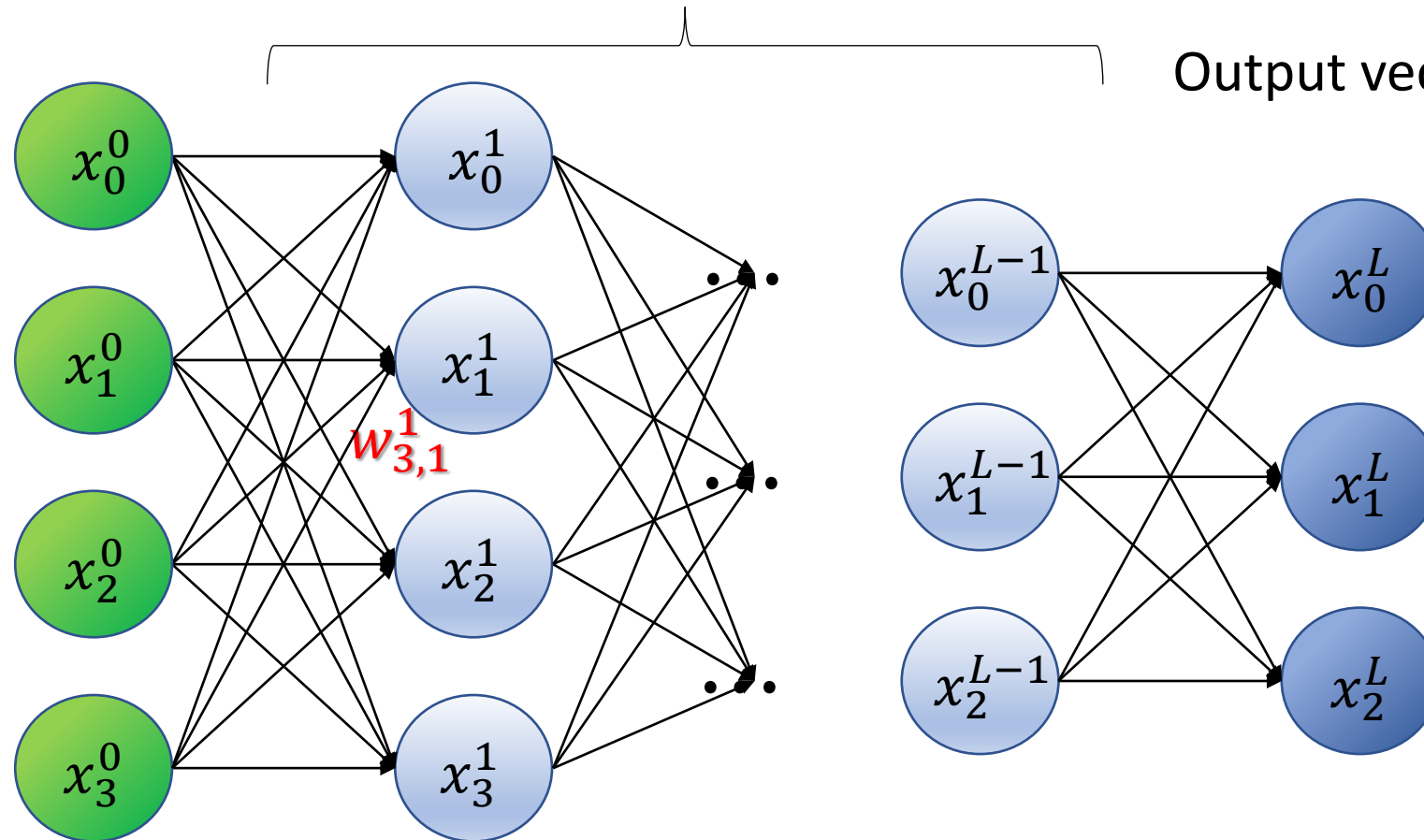
# Artificial Neural Networks

Vaguely inspired by biological neural networks

Input vector

Hidden layers

Output vector

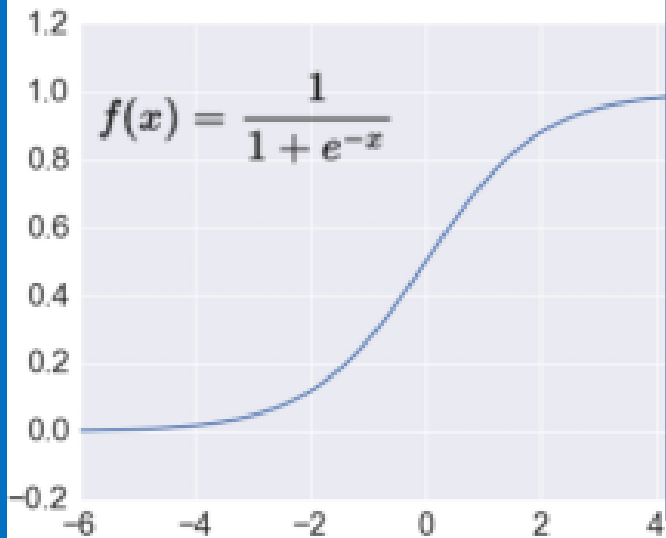


Q: What do you call a single layered net?

Q: Why?

Sigmoid

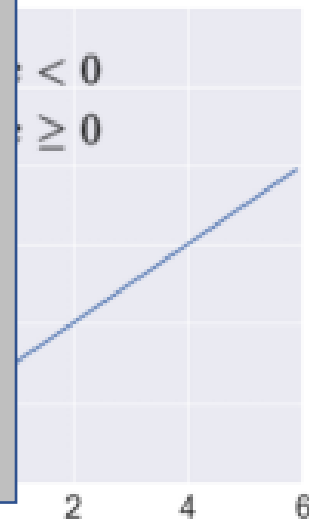
$$f(x) = \frac{1}{1 + e^{-x}}$$



Logistic Regression

Single-layer feed-forward neural network with sigmoid activation.

$x < 0$   
 $x \geq 0$



Linear

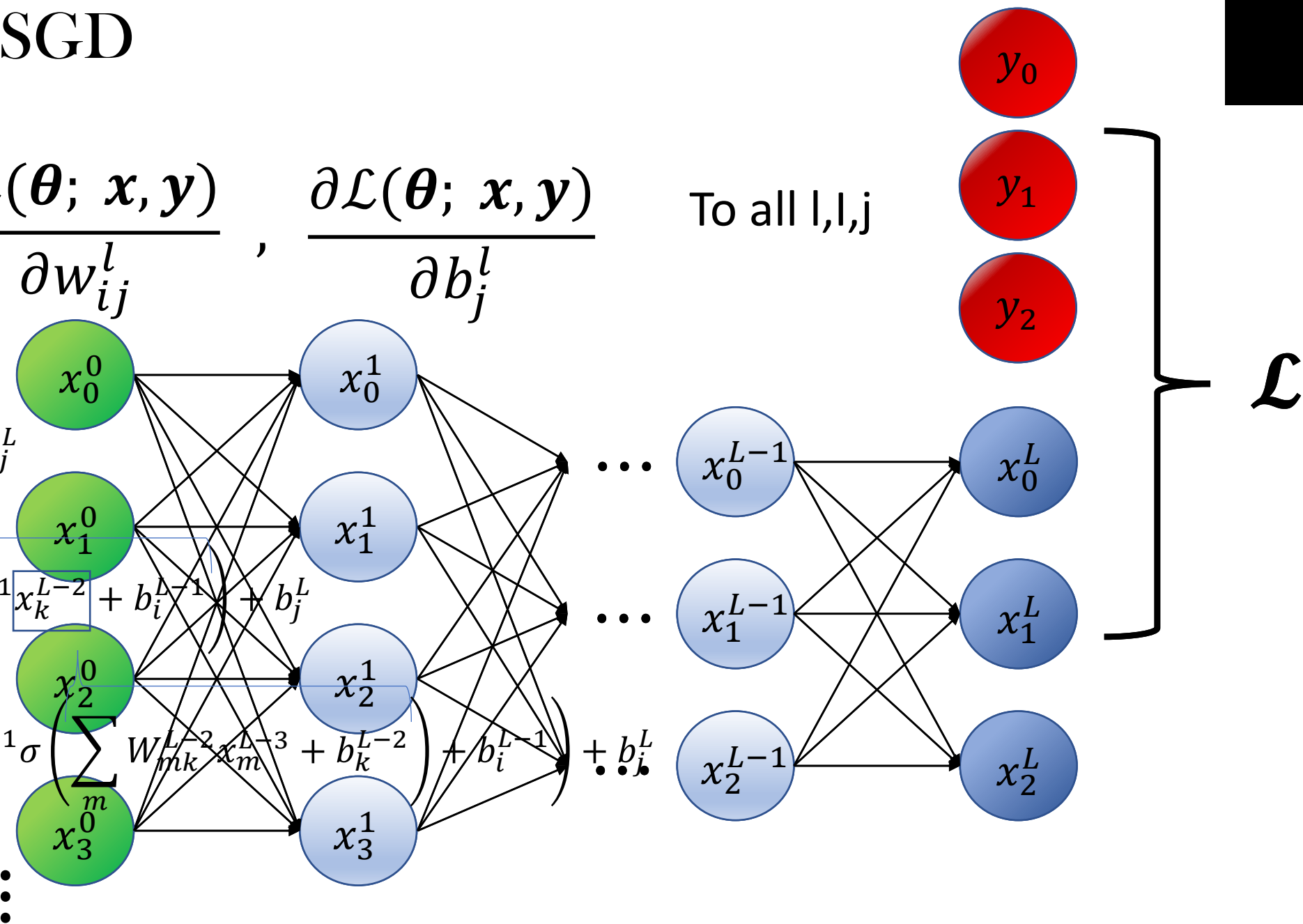
Feature transform

# Learning by SGD

We need  $\frac{\partial \mathcal{L}(\theta; \mathbf{x}, \mathbf{y})}{\partial w_{ij}^l}$ ,  $\frac{\partial \mathcal{L}(\theta; \mathbf{x}, \mathbf{y})}{\partial b_j^l}$

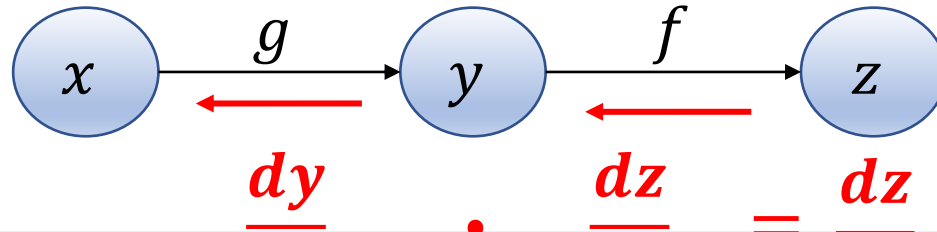
To all  $l, i, j$

$$\begin{aligned}
 x_j^L &= \sum_i W_{ij}^L x_i^{L-1} + b_j^L \\
 &= \sum_i W_{ij}^L \sigma \left( \sum_k W_{ki}^{L-1} x_k^{L-2} + b_i^{L-1} \right) + b_j^L \\
 &= \sum_i W_{ij}^L \sigma \left( \sum_k W_{ki}^{L-1} \sigma \left( \sum_m W_{mk}^{L-2} x_m^{L-3} + b_k^{L-2} \right) + b_i^{L-1} \right) + b_j^L
 \end{aligned}$$



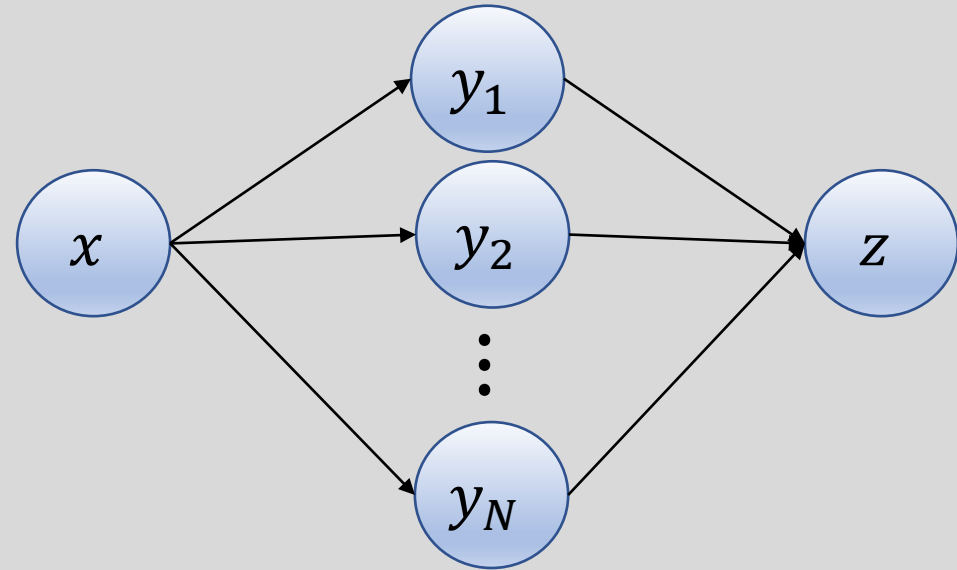
# Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$



Conclusion:

$$\frac{df(y_1, y_2 \dots y_N)}{dx} = \sum_n \frac{\partial f}{\partial y_n} \frac{dy_n}{dx}$$



$$= 6x^2 + 4e^{2x}$$

$$\frac{dz(y_1, y_2)}{dx} = 12x + 8e^{2x}$$

# Back Propagation - preliminaries

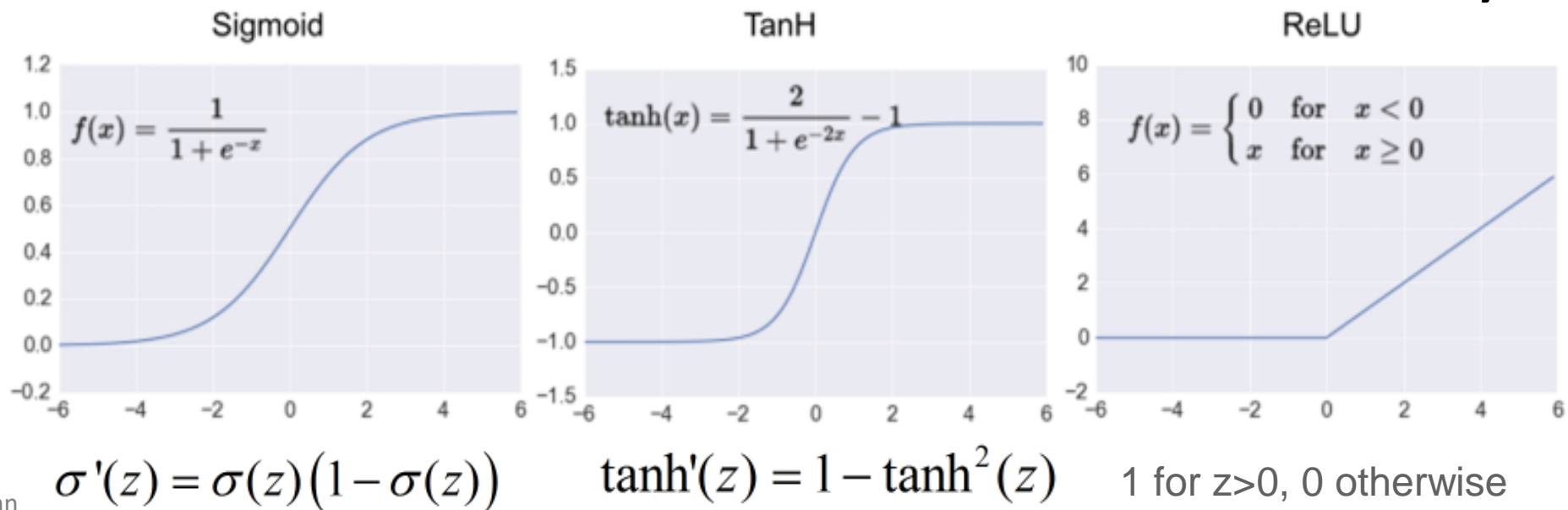
$$x_j^l = \sigma \left( \underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^l} = \frac{\partial \mathcal{L}}{\partial x_j^l} \cdot \frac{\partial x_j^l}{\partial w_{ij}^l}$$

$\triangleq g_j^l$   
Obtained by  
backprop

Easy!  
 $x_i^{l-1} \cdot \sigma'(z_j^l)$

## Derivatives of common activations are easy!



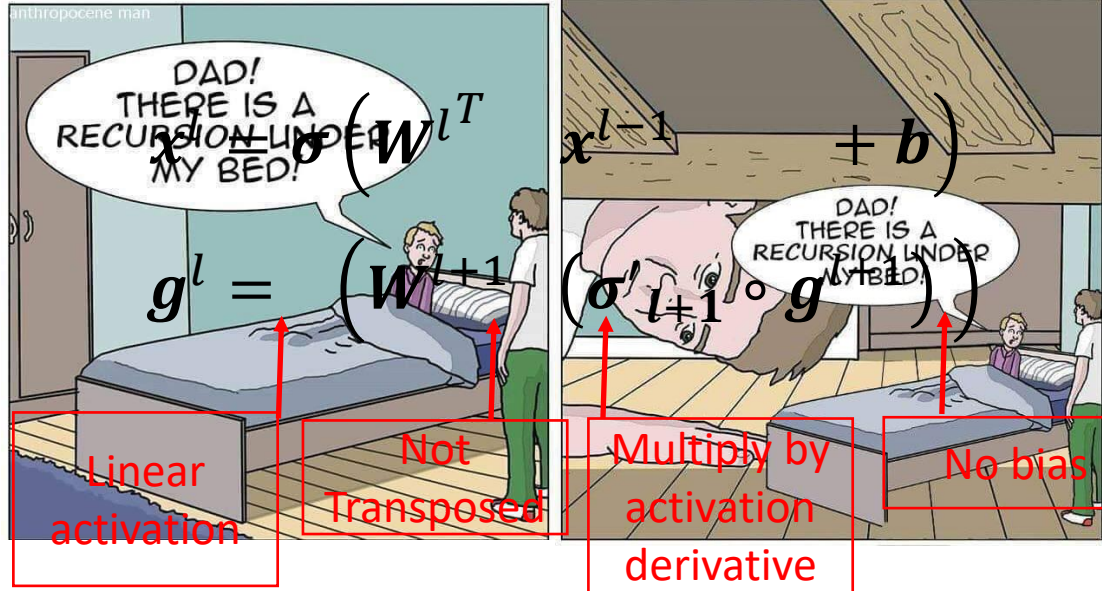


# Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})$$



Stopping criterion:

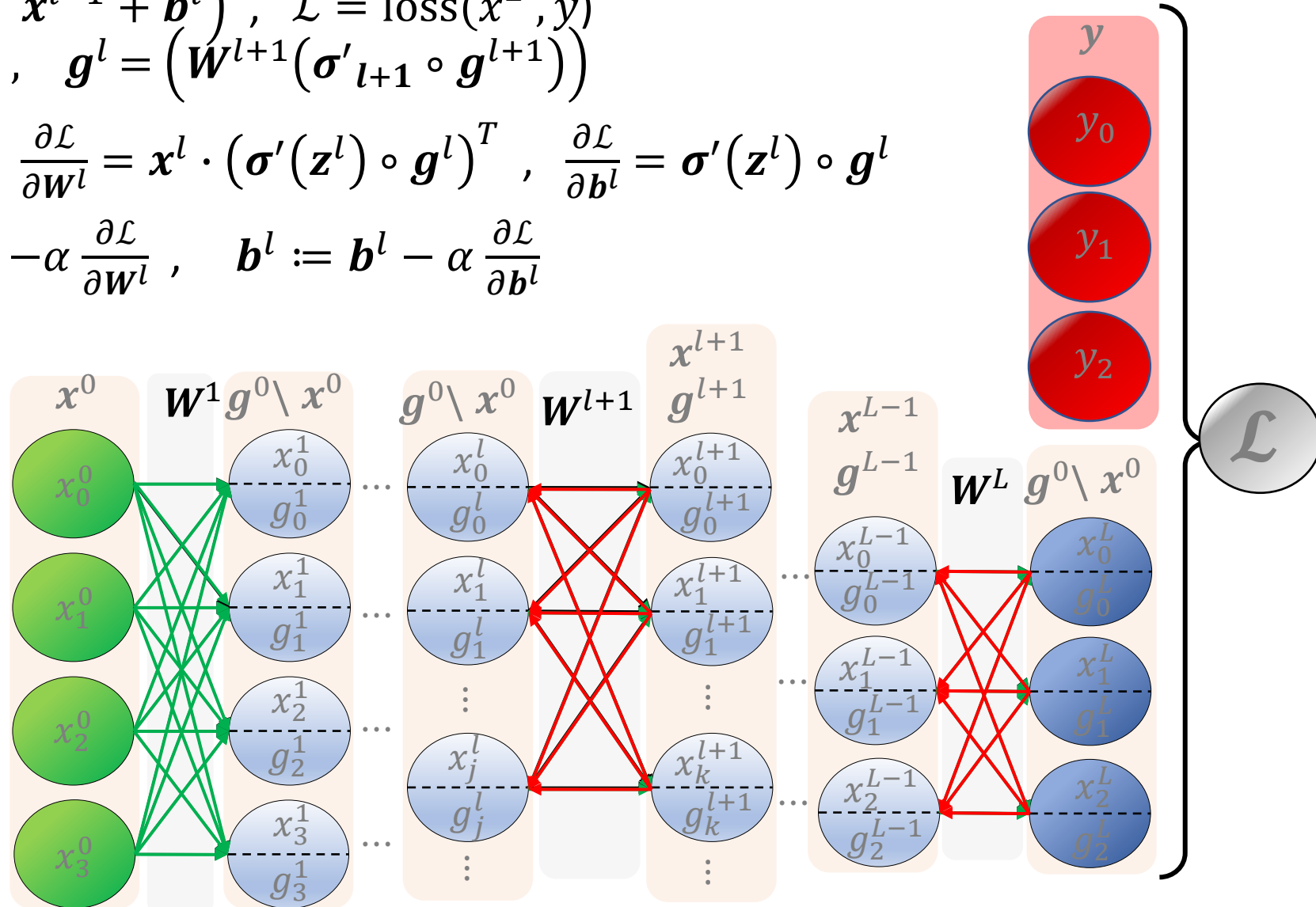
$$x_j^l = \sigma \left( \sum_i w_{ij}^l \cdot x_i^{l-1} + b_j \right)$$

$$g_j^l = \frac{\partial \mathcal{L}}{\partial x_j^l} \cdot \underbrace{\left( \sum_i w_{ij}^l \cdot x_i^{l-1} + b_j \right)}_{z_j^l}$$



- Initialize weights
- Repeat until convergence:
  1. Sample a batch from the data:  $\{(x_i, y_i) \dots\}$
  2. Forward pass:  $x^l = \sigma(W^{lT} x^{l-1} + b^l)$ ,  $\mathcal{L} = \text{loss}(x^L, y)$
  3. Backward pass:  $g^L = \frac{\partial \mathcal{L}}{\partial x^L}$ ,  $g^l = (W^{l+1}(\sigma'_{l+1} \circ g^{l+1}))$
  4. Calculate weights gradient:  $\frac{\partial \mathcal{L}}{\partial W^l} = x^l \cdot (\sigma'(z^l) \circ g^l)^T$ ,  $\frac{\partial \mathcal{L}}{\partial b^l} = \sigma'(z^l) \circ g^l$
  5. Update weights:  $W^l := W^l - \alpha \frac{\partial \mathcal{L}}{\partial W^l}$ ,  $b^l := b^l - \alpha \frac{\partial \mathcal{L}}{\partial b^l}$

# Back Propagation



# Let's get more generic



**Yann LeCun**

January 5, 2018 · 🌐

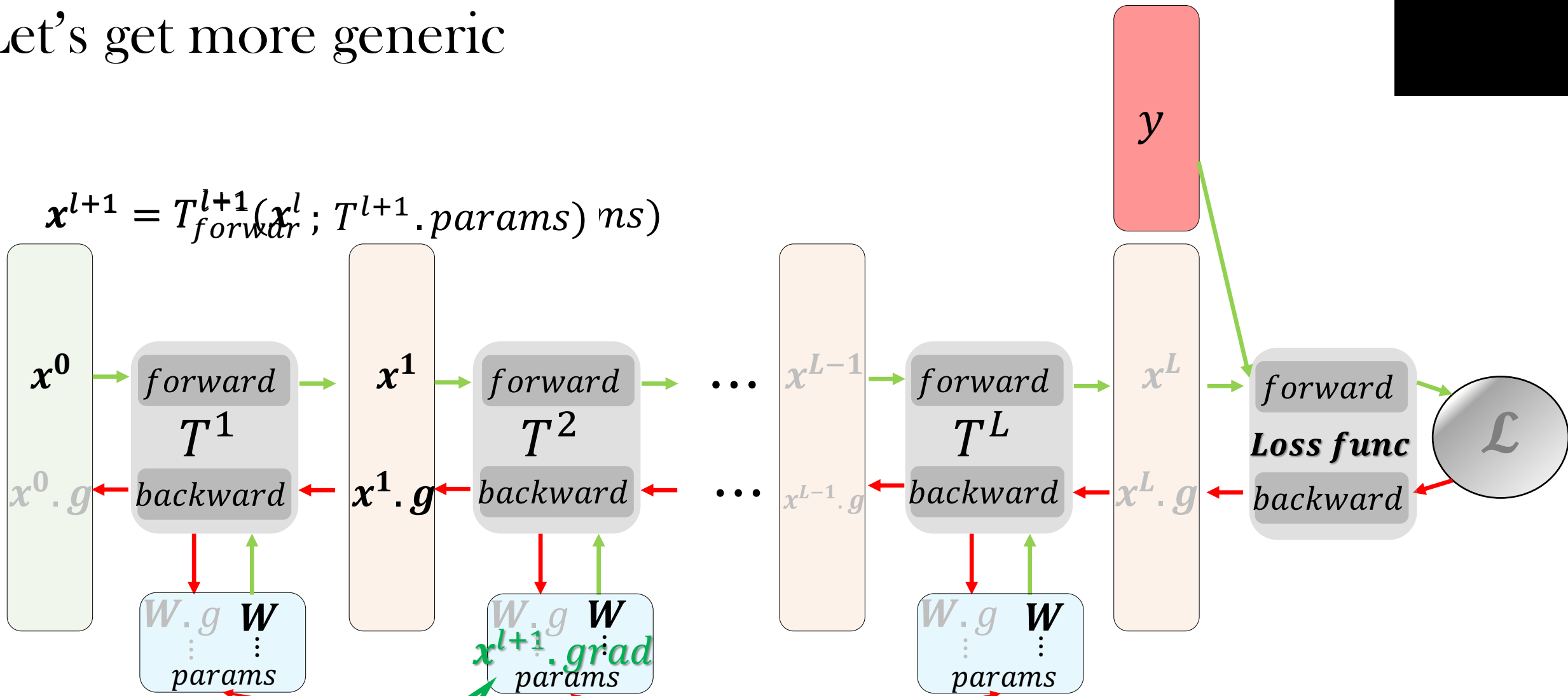


OK, Deep Learning has outlived its usefulness as a buzz-phrase. Deep Learning est mort. Vive **Differentiable Programming!**

Yeah, Differentiable Programming is little more than a rebranding of the modern collection Deep Learning techniques, the same way Deep Learning was a rebranding of the modern incarnations of neural nets with more than two layers.

But the important point is that people are now building a new kind of software by assembling networks of **parameterized functional blocks** and by training them from examples using some form of gradient-based optimization.

# Let's get more generic

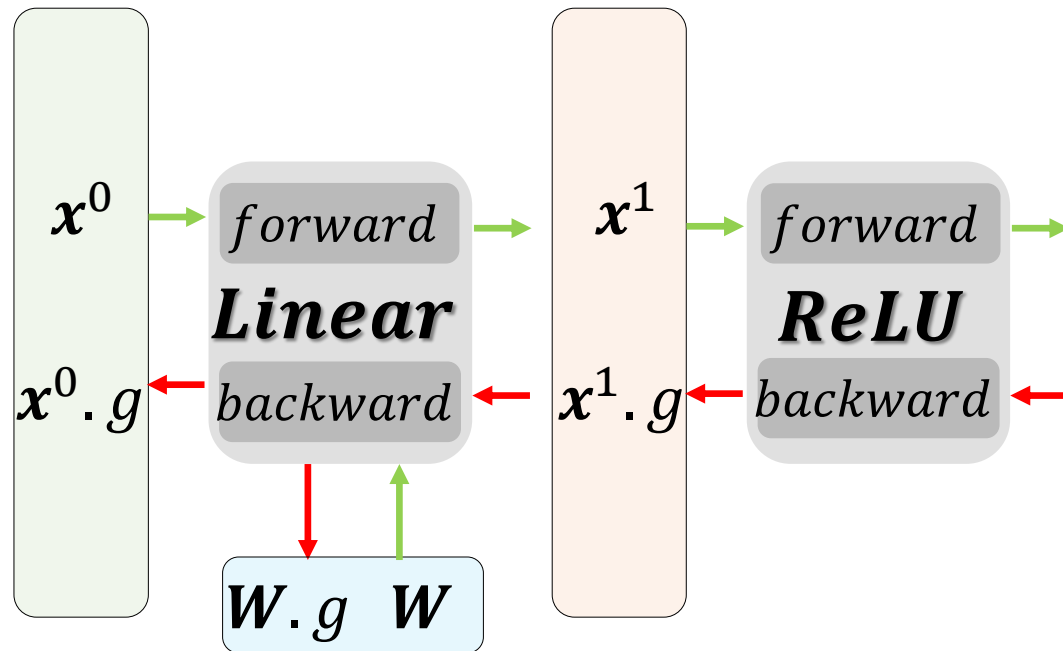


$$x^{l+1} = T_{forward}^{l+1}(x^l; T^{l+1}.params)$$

$$x^l.g = \frac{\partial \mathcal{L}}{\partial x^l} = \frac{\partial \mathcal{L}}{\partial x^{l+1}} \cdot \frac{\partial x^{l+1}}{\partial x^l} = T_{backward}^{l+1}(x^{l+1}.grad; T^{l+1}.params, x^l)$$

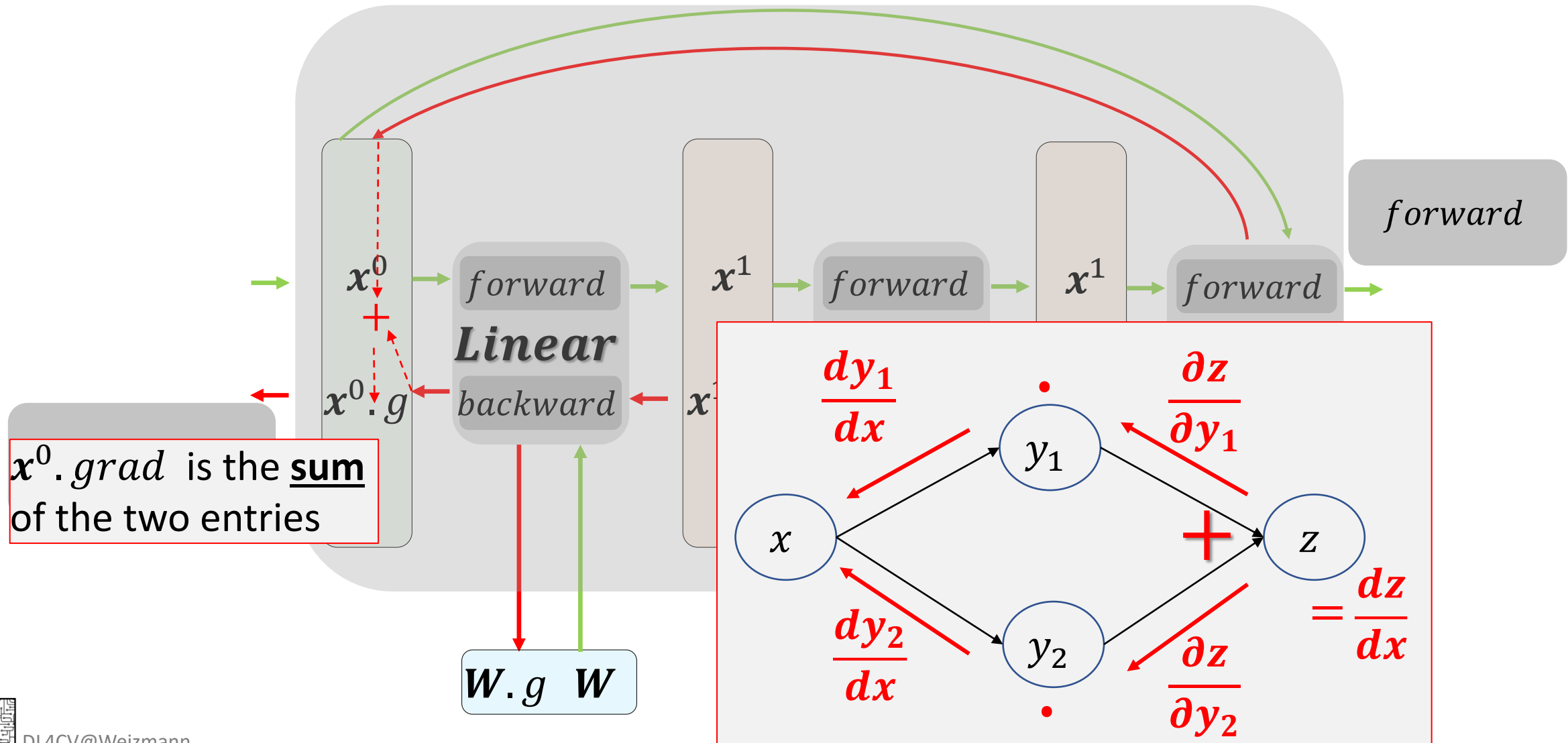
Use to update weights.  
 $W := W - \alpha W.g$

# Example: Standard layer



BTW : You can backprop any DAG!

BTW2: Layers (NN modules) can be nested!



# Yes you should understand

Be creative,  
but always watch your back(prop)!



<http://playground.tensorflow.org>





This week's tutorial:



Dolev Ofri

# Intro to PyTorch

Next week's lecture:

(Me  
Again ☹ )

# Convolutional Neural Networks

